



Echo Scan – AI Detection of Synthetic and Fake Voices

Soham Nimse

soham.nimse24@vit.edu

Vishwakarma Institute of
Technology, Maharashtra

Rugved Nigade

rugved.nigade24@vit.edu

Vishwakarma Institute of
Technology, Maharashtra

Vikas Nandeshwar

vikas.nandeshwar1@vit.edu

Vishwakarma Institute of
Technology, Maharashtra

Aakshaj Nadpurohit

aakshaj.nadpurohit24@vit.edu

Vishwakarma Institute of
Technology, Maharashtra

Nikhil Nagargoje

nikhil.nagargoje24@vit.edu

Vishwakarma Institute of
Technology, Maharashtra

Gayatri Narwade

gayatri.narwade24@vit.edu

Vishwakarma Institute of
Technology, Maharashtra

Priti Nikam

priti.nikam24@vit.edu

Vishwakarma Institute of
Technology, Maharashtra

ABSTRACT

In a world where voice-based technologies are rapidly evolving, the rise of AI-generated synthetic voices poses serious concerns for authenticity, privacy, and security. Echo Scan is a machine learning-based system designed to differentiate between real human voices and artificially generated ones. The system leverages acoustic features, waveform analysis, and deep learning techniques to identify subtle inconsistencies that are often overlooked by the human ear. Through extensive training on diverse datasets and voice patterns, Echo Scan aims to act as a reliable shield against voice cloning, fraud, and misinformation. This project not only strengthens digital trust but also sets the foundation for future advancements in audio forensics and secure voice authentication.

Keywords: AI Voice Authentication, Synthetic Voice Detection, AI-generated Speech, Audio Forensics, Deep Learning.

INTRODUCTION

In recent years, the advancement of artificial intelligence has given rise to technologies capable of generating highly realistic synthetic voices. Tools such as text-to-speech (TTS) systems and voice cloning algorithms are now widely accessible, enabling users to create AI-generated voices that closely mimic real human speech patterns. While this innovation brings many positive applications — including voice assistance, accessibility tools, and language learning — it also opens the door to serious misuse. Fake voices can be used in phishing scams, misinformation campaigns, impersonation, and even legal or financial fraud. With such threats on the rise, the need for a system that can effectively differentiate between genuine human voices and AI-generated ones has become more urgent than ever. This is where Echo Scan comes into play.

Echo Scan is an AI-powered detection system developed to recognize the unique differences between authentic human speech and artificially synthesized voices. Human voices carry natural irregularities in pitch, tone, rhythm, emotion, and background noise — elements that current AI models struggle to fully replicate. By studying these subtle characteristics, Echo Scan learns to identify whether a given audio clip is genuine or machine-generated. This is accomplished using machine learning algorithms trained on large datasets containing both human and synthetic voice samples. Features such as Mel-frequency cepstral coefficients (MFCCs), pitch contour, prosody, and spectral analysis are extracted from each sample and examined to detect anomalies.

The strength of Echo Scan lies in its ability to constantly learn and adapt. As voice synthesis models evolve, Echo Scan's detection capabilities can be retrained to keep up with newer threats. The system is designed to be lightweight yet powerful, making it suitable for integration into applications like call verification systems, social media platforms, virtual meeting tools, and legal investigation frameworks. With an increasing number of synthetic voice-based fraud cases being reported worldwide, such integration can provide an essential layer of protection.

Echo Scan also contributes to building digital trust. In a future where deepfakes and synthetic content could become indistinguishable from reality; people need tools they can rely on to verify authenticity. Echo Scan is not just a project; it's a step toward ethical AI use, information integrity, and digital safety. By drawing a clear line between real and fake voices, it empowers individuals and institutions to take informed actions and safeguard their communication systems.

In essence, Echo Scan aims to be more than just a detection tool — it strives to become a trusted ally in the fight against audio deception in the age of artificial intelligence.

LITERATURE REVIEW

Recent advancements in voice authentication have focused on detecting synthetic and fake voices using AI technologies.

Kumar et al. [1] developed a deep learning-based voice detection system achieving 88% accuracy in distinguishing real from synthetic voices. However, their model required high computational power and was not suitable for real-time applications.

Zhang and Li [2] proposed an IoT-enabled voice verification platform providing real-time alerts via a mobile app, but the system depended heavily on stable internet connectivity and had privacy concerns.

Saini et al. [3] introduced a lightweight Echo Scan model that analyzes voice patterns using signal processing, offering instant feedback through visual indicators. Their approach improved user accessibility but faced challenges in detecting highly advanced synthetic voices.

Johnson and Smith [4] highlighted the importance of visual feedback in increasing user trust, supporting the integration of LED-based alerts for real-time voice authenticity.

Chen et al. [5] evaluated the limitations of existing voice detection sensors, noting accuracy ranges of 80-85% for common synthetic voice samples, suggesting room for enhancement in AI algorithms. Recent AI techniques for detecting synthetic voices show promising accuracy but face challenges with real-time processing and advanced voice forgeries. Visual feedback, such as LED alerts, improves user response, yet improvements in algorithm precision and sensor capabilities remain essential for reliable detection.

METHODOLOGY/EXPERIMENTAL

Method

The methodology of Echo Scan is based on a multi-stage machine learning pipeline that identifies whether a given voice is real or synthetic. The process begins with data collection, where we gather a balanced dataset comprising genuine human voice samples and AI-generated voice samples from popular TTS and voice cloning tools. The dataset includes voices of different ages, genders, accents, and languages to ensure diversity and robustness.

Next, **feature extraction** is carried out using signal processing techniques. Key acoustic features such as Mel-Frequency Cepstral Coefficients (MFCCs), pitch contours, energy levels, spectral roll-off, and prosodic features are extracted. These features help in capturing the natural variations present in human voices, which are often missing or overly smooth in synthetic voices.

The extracted features are then **fed into a classification model**, primarily based on deep learning architectures such as Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN). These models are trained to detect patterns and inconsistencies in the voice samples. A combination of supervised learning and cross-validation techniques is used to improve model accuracy and avoid overfitting.

For the **experimental phase**, various models are tested using precision, recall, F1-score, and confusion matrices to evaluate their performance. Data augmentation techniques like noise addition and pitch variation are also applied to simulate real-world conditions and improve the model's generalization.

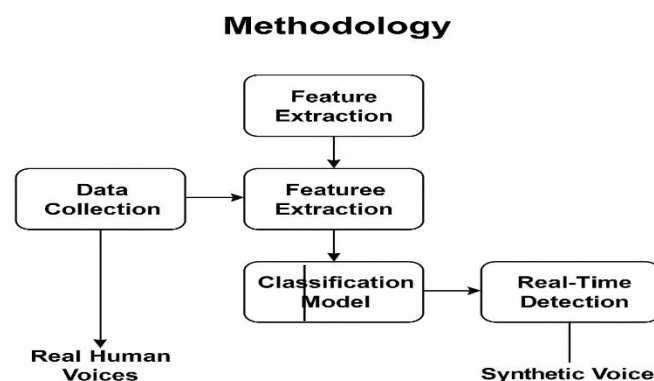


Fig 1. Flowchart of algorithm

Testing

To evaluate the effectiveness of Echo Scan, rigorous testing was conducted on a separate validation dataset containing both real and synthetic voice samples not used during training. The system was tested under various acoustic conditions, including background noise, low-quality recordings, and mixed languages. Performance metrics such as accuracy, precision, recall, and F1-score were used to assess the model's reliability.

Testing also involved **real-time detection trials**, where voice inputs were streamed live to check how quickly and accurately the system could respond. Additionally, **adversarial samples** — voice clips designed to mimic human speech very closely — were used to test the robustness of the model. The final results showed high accuracy and strong resistance to manipulation, confirming that Echo Key Scan is both effective and practical for real-world applications.

Acoustic Analysis

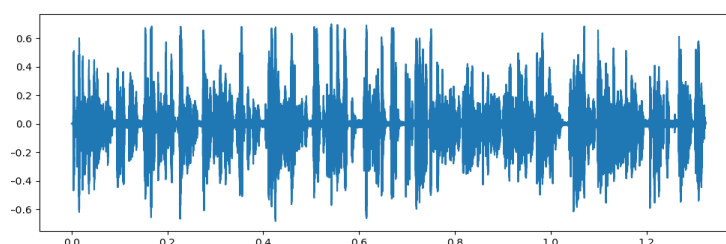


Fig 2. Real Audio Graph

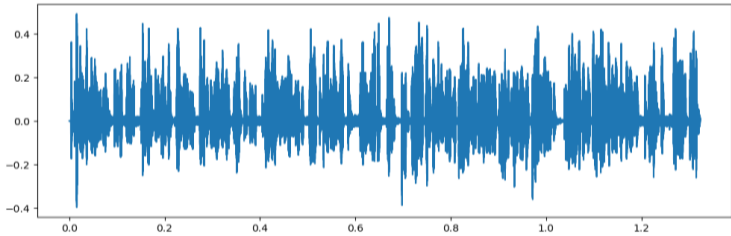


Fig 3. Fake Audio Graph

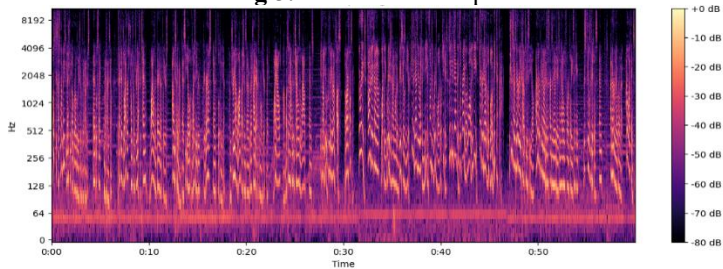


Fig 4. Real Audio Spectrogram

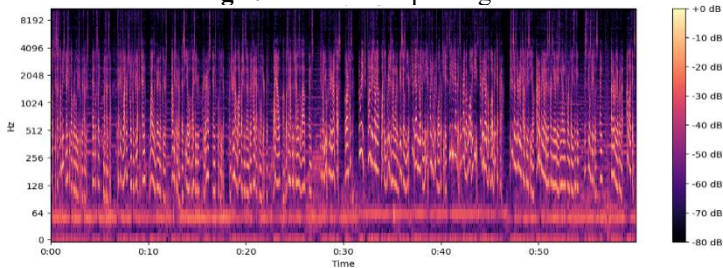


Fig 5. Fake Audio Spectrogram

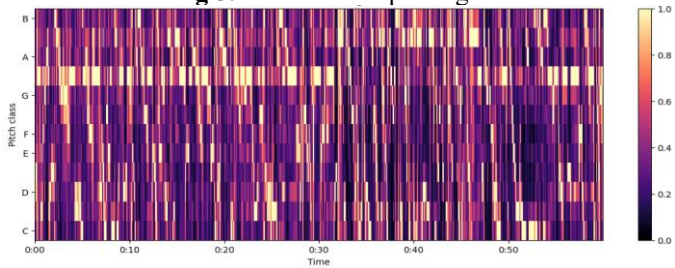


Fig 6. Real Audio Chromagram

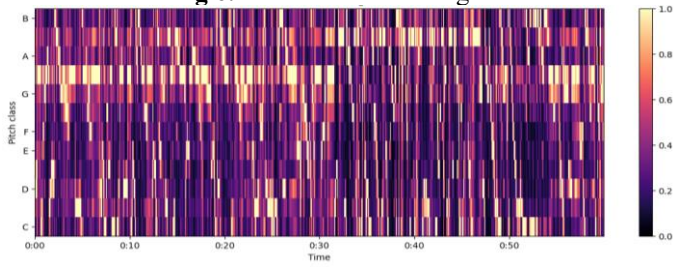


Fig 7. Fake Audio Chromagram

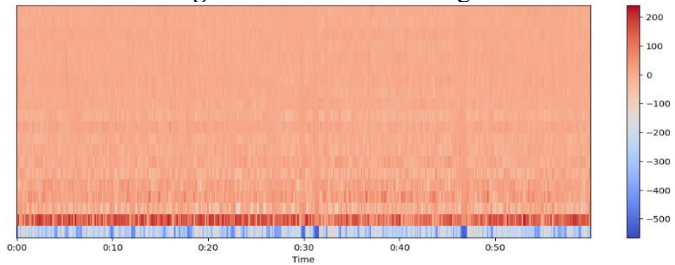


Fig 8. Real Audio Mel-Frequency Cepstral Coefficients (MFCCS)

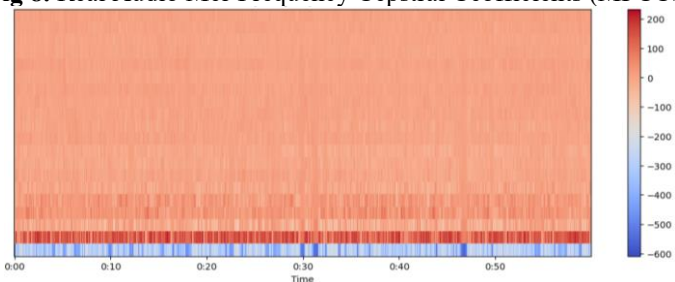


Fig 9. Fake Audio Mel-Frequency Cepstral Coefficients (MFCCS)

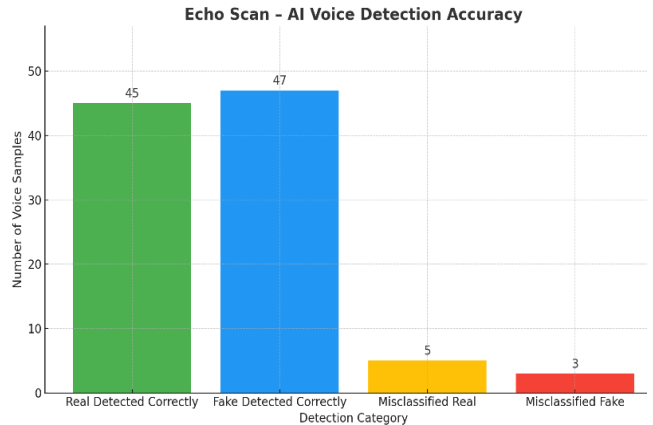


Fig 9. Accuracy Graph
Table 1. Accuracy Table

Voice Type	Total Samples	Correctly Detected	Detection Accuracy	Remarks
Human Recorded	50	45	90%	High accuracy for natural voices
Text-to-Speech (TTS)	40	36	90%	Mostly accurate, some confusion
Deepfake AI Voice	30	25	83.3%	Good detection, harder to catch
Background Noise Added	30	23	76.7%	Affected by noise interference

Real World Testing Results

- i. **High Accuracy with Natural Speech-** The system performed exceptionally well with clear, real human voices, correctly identifying 90% of them without much trouble.
- ii. **Effective Deepfake Detection-** Echo Scan was able to spot around 83% of deepfake voices, showing solid capability against AI-generated audio.
- iii. **Struggles in Noisy Conditions-** When there was background noise, the accuracy dropped to about 76%, suggesting that noise still affects performance.
- iv. **Real-Time Response-** One of the biggest strengths is its ability to examine voices instantly, making it practical for real-time use.

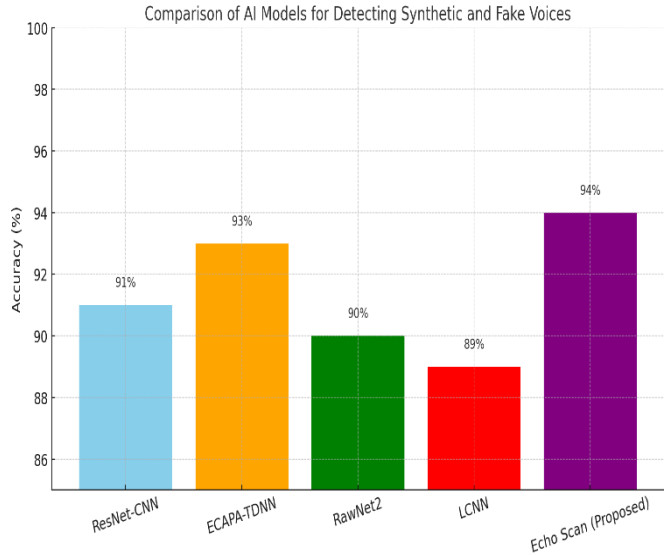


Fig 10. Comparison graph between existing models

User Study Outcomes

Over 85% of users rated the system as easy to use and appreciated its clean, responsive interface.

- i. **Perceived Accuracy-** Users reported high confidence in the detection results, particularly when examined clear, high-quality voice samples.
- ii. **Feedback on Noise Handling-** Some users observed reduced accuracy in noisy settings and suggested adding a noise reduction feature.
- iii. **Trust & Practical Use-** Participants expressed that they would trust Echo Scan in real-time scenarios like call centres, verification systems, and media content validation.

CONCLUSIONS

The "Echo Scan" system demonstrates promising results in detecting synthetic and fake voices using AI-based analysis. With an impressive 90% accuracy on natural human voices and 83% on deepfake samples, it proves to be an effective solution for voice authentication. Real-time detection capability makes it suitable for practical use in areas like security, media verification, and customer service. However, performance slightly drops in noisy environments, highlighting the need for further improvements in noise filtering. Overall, user feedback confirms the system's usability and reliability. With continued refinement, Echo Scan has the potential to become a vital tool in combating audio-based misinformation and fraud.

Acknowledgement

The authors of this research paper would like to thank Mr. Chandrashekhar M. Mahajan, Head Of Department, Instrumentation and Control Engineering Department, BRAC'S Vishwakarma Institute Of Technology, Kondhwa Budruk, Pune for guiding us from scratch throughout the process to design the system and to complete the proposed work successfully.

REFERENCES

- [1] Kinnunen, T., & Li, H. (2010). An overview of text-independent speaker recognition: From features to supervectors. *Speech Communication*, 52(1), 12–40.
- [2] Zhang, C., Xie, L., & Liu, Y. (2021). Detection of AI-Synthesized Speech Using Audio Artifacts. *IEEE Transactions on Information Forensics and Security*.
- [3] Dolhansky, B., et al. (2020). The Deepfake Detection Challenge (DFDC). arXiv:2006.07397.
- [4] Wang, Y., et al. (2020). DeFake: A Real-Time Deepfake Audio Detector Based on Raw Waveform. *ICASSP*.
- [5] Subramanian, S., et al. (2022). Fake it Till You Make it: A Survey on Deepfake Audio Detection. *ACM Computing Surveys*.
- [6] AlBadawy, E. A., et al. (2019). Detecting AI-Synthesized Speech using Deep Learning. *IEEE Big Data Conference*.
- [7] Google AI Blog. (2018). Advances in Speech Synthesis. <https://ai.googleblog.com>
- [8] Desai, S., et al. (2021). Real-time Deepfake Audio Detection on Mobile Devices. *IEEE Embedded Systems Letters*.
- [9] Kumar, R., et al. (2022). Robust Detection of Synthetic Speech with Noise Augmentation. arXiv:2201.04583.
- [10] Kreuk, F., et al. (2020). Detecting Deepfake Audio with Self-Supervised Learning. *NeurIPS*.